1 **Structural and functional analysis of genes with**

2 **potential involvement in resistance to coffee leaf rust: a**

3 **functional marker based approach**

4

5 **Geleta Dugassa Barka[1,4], Eveline Teixeira Caixeta[1,2, *], Sávio Siqueira Ferreira[,1], Laércio**

6  **Zambolim[3]**

7

8 [1]Laboratório de Biotecnologia do Cafeeiro (BIOCAFÉ), BIOAGRO, Universidade Federal de Viçosa

9 (UFV), 36570-000, Viçosa, MG, Brazil

10 [2]Embrapa Café, Empresa Brasileira de Pesquisa Agropecuária, 70770-901, Brasília, DF, Brazil

11 [3]Departamento de Fitopatologia, Universidade Federal de Viçosa (UFV), 36570-000, Viçosa, MG, Brazil

12 [4]Applied Biology Department, Adama Science and Technology University (ASTU), 1888, Adama,

13 Oromia, Ethiopia

14

15 *Corresponding author

16 E-mail: eveline.caixeta@embrapa.br (ETC)

17

18

19

20

21

22

23

# Abstract

Physiology-based differentiation of $S_H$ genes and *Hemileia vastatrix* races is the principal method employed for the characterization of coffee leaf rust resistance. Based on the gene-for-gene theory, nine major rust resistance genes ($S_H$1-9) have been proposed. However, these genes have not been characterized at the molecular level. Consequently, the lack of molecular data regarding rust resistance genes or candidates is a major bottleneck in coffee breeding. To address this issue, we screened a BAC library with resistance gene analogs (RGAs), identified RGAs, characterized and explored for any $S_H$ related candidate genes. Herein, we report the identification and characterization of a gene (gene 11), which shares conserved sequences with other $S_H$ genes and displays a characteristic polymorphic allele conferring different resistance phenotypes. Furthermore, comparative analysis of the two RGAs belonging to CC-NBS-LRR revealed more intense diversifying selection in tomato and grape genomes than in coffee. For the first time, the present study has unveiled novel insights into the molecular nature of the $S_H$ genes, thereby opening new avenues for coffee rust resistance molecular breeding. The characterized candidate RGA is of particular importance for further biological function analysis in coffee.

# Introduction

Coffee is one of the most valuable cash crops in many developing economies as it provides employment opportunities in cultivation, processing and marketing activities, thereby sustaining the livelihoods of millions around the world [1]. *H. vastatrix*, the causative agent of coffee leaf rust, accounts for one of the major threats to coffee production in almost every coffee producing region. Despite the release of some resistant coffee cultivars in recent years, coffee rust continues to adversely affect coffee production and undermines the incomes of many households [2]. To date, at least 49 characterized physiological races of *H. vastatrix* have been reported [2,3]. The consistent emergence of new races and the sporadic outbreaks of this disease have imposed challenges in resistance breeding. The most pressing concern is, however, the breakdown of resistance genes leading to disease susceptibility of cultivars that were once validated as superior genetic material for resistance breeding [4].

50    The molecular profiles of coffee genes involved in different metabolic pathways, their evolution

51    and annotation have been unveiled with the complete sequencing of *C. canephora* genome [5]. Given that

52    *C. canephora* contributes to half of the Arabica coffee genome, being a natural hybrid of *C. canephora*

53    and *C. eugenioides*, open access to its genome has provided valuable insights into the genome of *C.*

54    *arabica* during the past five years. The discovery and successful introgression of $S_H3$ resistance gene

55    locus to cultivated Arabica coffee from *C. liberica* was another landmark often considered as one of the

56    greatest milestones in the development of coffee rust resistance [6]. Since then, molecular and physical

57    mapping has enabled the sequencing and annotation of $S_H3$ region, resulting in the discovery of multiple

58    resistance (R) genes [6,7]. Dominantly inherited, the largest class of R-genes encode nucleotide-binding

59    site leucine-rich-repeat (NBS-LRR) proteins that directly recognize the corresponding virulence (v)

60    protein of the pathogen or its effects [8,9]. These genes are believed to contain several hundred gene

61    families, which are unevenly distributed in the genomes of different plant species [10]. Intracellular

62    signaling domain, similar to Drosophila toll/mammalian interleukin-1 receptor (TNL, Toll-NBS-LRR)

63    and the coiled-coil (CNL, CC-NBS-LRR), are the two major N-terminal amino acid sequences preceding

64    the NBS domain involved in specific signal transduction [8,11]. The other N-terminal domain linked to

65    LRR includes leucine-zipper (a transmembrane protein, TM), protein kinase (PK) and WRKY TIR

66    proteins [12]. These domains are predominantly involved in resistance signal transduction via

67    conformational changes [13]. On the carboxyl-terminal region is the LRR, mediating specific protein-

68    protein interaction to recognize pathogen effectors [14,15]. Although these domains are few in number,

69    nucleotide polymorphism and variability of the LRR region are responsible for the perception of a

70    specific pathogen effector [9,16]. Inter and intraspecific extreme variabilities of NBS-LRR have been

71    attributed to gene duplication, unequal crossing over, recombination, deletion, point mutation and

72    selection pressure due to continuous response to diverse pathogen races [6].

73    The readily available Arabica coffee BAC libraries constructed from disease resistant genotypes

74    at different laboratories have accelerated studies involving resistance gene cloning [17,18]. Furthermore,

75    the application of arbitrary DNA-based and functional (gene) markers in gene cloning has benefitted crop

76    improvement, either through map-based cloning using the former or direct gene cloning using the latter or

77    both [19]. Direct cloning of the gene of interest over map-based gene cloning is appealing as this method

78    is more precise and straightforward for gene characterization.

79    In coffee, the origin and organization of disease resistance genes have begun to emerge in recent

80    years as part of an effort to understand the role of major rust resistance genes. One such endeavor was the

81    assembly of R genes spanning the $S_H3$ locus with the objective of tracing the evolution and diversity of

82    LRR domains in three coffee species [6]. Despite the partial sequencing and annotation of several disease

83    resistance genes in Arabica coffee [20], completely sequenced and characterized candidate genes are not

84    yet readily available. Resistance to rust is conferred by nine major genes ($S_H$1-9) and the corresponding

85    $v_{1-9}$ pathogen factors are known for long in the coffee rust pathosystem [3,21]. Nonetheless, molecular

86    and functional characterization of any of the $S_H$ genes and the associated regulatory elements is entirely

87    obscure, yet holds immense potential in changing the perspective of rust resistance breeding. Likewise,

88    the use of functional markers that serve as a direct rust resistance screening tool amongst the host-

89    differential coffee clones is important but is barely addressed. The lack of a typical candidate rust

90    resistance gene is one of the bottlenecks in coffee breeding. A resistance gene analog (RGA) marker,

91    CARF005, was previously confirmed to share disease resistance ORF region in coffee [20,22]. This

92    polymorphic RGA marker encodes the disease resistance protein domain NB-ARC (nucleotide binding

93    site-ARC: ARC for APAF-1, R protein and CED-4 [23], exclusive in coffee cultivars resistant to *H.*

94    *vastatrix* [22]. The complete sequencing and molecular characterization would help identify candidate

95    disease resistance genes. The state-of-the-art bioinformatics analysis, availability of differential coffee

96    clones with specific $S_H$ genes, structural and functional analysis of conserved domains and associated

97    motifs of candidate RGAs belonging to the $S_H$ gene series could greatly advance coffee rust resistance

98    breeding. Therefore, the objectives of this study were to trace the origin of resistance gene analogs

99    (RGAs) involved in coffee rust resistance and perform comparative molecular characterization of selected

100    candidate gene to determine whether it belongs to the $S_H$ gene series. We also investigated if any of the

101    RGAs were activated during incompatible interaction between C. *Arabica* and *H. vastatrix*..

# Materials and methods

## Plant materials

104    Twenty-one differential coffee clones containing at least one of the coffee rust resistance genes

105    ($S_H$1-9) and three genotypes susceptible to all the virulence factors (v1-9) of *H. vastatrix* were used in the

106    CARF005 screening. The differential clones were initially characterized by CIFC (Centro de Investigação

4

107    das Ferrugens do Cafeeiro, Portugal) for the identification of the different physiological races of *H.*

108    *vastatrix*. All clones were vegetatively propagated at the Plant Pathology Department greenhouse of the

109    Universidade Federal de Viçosa (Brazil). Genomic DNA was extracted from a young, second pair of

110    leaves following [24]. DNA integrity was checked by electrophoresing in 1% gel and visualized after

111    staining with ethidium bromide (0.5 µg/ml) and Nanodrop (NanoDrop Technologies, Wilmington, USA).

112    DNA was stored at -20 °C until further use. RNA-Seq libraries (hereafter, referred to as transcriptome)

113    were constructed at 12 and 24 h after infection (hai) during the *C. arabica* CIFC 832/2 –*H. vastatrix* (race

114    XXXIII) incompatible interaction [25] and were used as the reference in the search for novel candidate

115    resistance genes.

116

## PCR conditions

118    A Sigma made (Sigma-Aldrich, Belo Horizonte, Brazil) disease RGA primer pair, CARF005, (F:

119    5'-GGACATCAACACCAACCTC-3' and R: 5'-ATCCCTACCATCCACTTCAAC-3') [26]  was used to

120    screen the differential host clones. PCR reagents were 1x buffer, 0.2 mM dNTPs, 0.2 µM primers, 1 mM

121    $MgCl_2$, 0.8 units of Taq polymerase (Invitrogen, Carlsbad, USA) to which 5 ng gDNA was added to form

122    a reaction volume of 20 µl. PCR cycling parameters were as follows: DNA denaturation at 95°C for 5

123    min followed by 35 cycles of 94°C for 30 s, 60°C for 30 s and 72°C for 1 min, followed by an extension

124    step at 72°C for 10 min. PCR products were screened for target inserts by electrophoresing in 1%

125    UltraPure$^{TM}$ agarose (Invitrogen) and visualized after staining with ethidium bromide (0.5 µg/ml). All

126    PCR and gel electrophoresis conditions were maintained consistently throughout the study unless stated

127    otherwise.

128

## Screening of BAC clone

130    BAC library comprising 56,832 clones, constructed using renowned rust resistant Hibrido de

131    Timor clone CIFC 832/2 [18] was used as the target source for RGA (CARF005). These clones were

132    replicated in 384-well titer plates using a plate replicator sterilized with 10% $H_2O_2$ for 2 min, rinsed in

133    sterile water for 10 seconds and soaked in 70% ethanol in laminar airflow cabinet. After the alcohol

134    evaporated (3-5 min), the old cultures were copied to a new 384-well titer plate containing 70 µl fresh LB

135    media (supplemented with 12.5 µgml$^{-1}$ chloramphenicol) in each well. Culture multiplication was

5

136    achieved by incubating the plates at 37 °C for 18 h and shaking at 180 rpm. The identification of clones

137    using the CARF005 insert was performed by grouping and subsequent group decomposition of the 384

138    clones until a single clone was identified as outlined in S1 Fig. BAC DNA was extracted using the

139    centrifugation protocol of Wizard® SV Plus Minipreps DNA Purification System (Promega, Fitchburg,

140    USA).

141

## Sequencing and contig assembly

143    The single BAC clone isolated using the CARF005 fragment was sequenced using the Illumina

144    HiSeq2000/2500 100PE (paired-end reads) platform at Macrogen (Seoul, South Korea). Paired-end

145    sequence processing and contig assembly were done using SPAdes software [27]. Contigs that matched

146    bacterial genome *(E. coli)* and sequences of the flanking vector (pCC1BAC$^{TM}$) were excluded prior to any

147    downstream sequence processing. The assembled BAC contigs were used to map against a transcriptome

148    constructed from coffee genes that were activated in response to *H. vastatrix* infection by Tophat 2 [28]

149    and to locate the contig region with active gene expression.

150

## Gene prediction and annotation

152    Contigs with $\geq$ 200 bp size and sharing $\geq$ 90% identity with *C. canephora* were subjected to

153    Augustus gene prediction [29]. Among the available genomes in the Augustus dataset, *Solanum*

154    *lycopersicum* was used as a reference genome, as they shared common gene repertoires and had similar

155    genome sizes [30]. The predicted ORFs were annotated using different online annotation tools. First,

156    BLASTp (NCBI) was used to detect the conserved domains and retrieve their description, followed by the

157    use of Predict Protein Server tool [31] molecular analysis and associated GO search. Protein 3D structure

158    and nucleotide (ATP/ADP/GTP/GDP) binding sites were predicted using I-TASSER suite online tool

159    [32]. As an annotation complement, the predicted ORFs were queried against the database to check for

160    coding sequences (CDS) of *S. lycopersicum* (Sol Genomics Network: https://solgenomics.net/tools/blast/)

161    and *V. vinifera* (Phytozome 11: https://phytozome.jgi.doe.gov) genomes.

162

## Sequence alignment and comparative analysis

6

164    Genes encoding resistance proteins were mapped to the *C. canephora* genome [5] to trace their

165    probable origin and organization. BLASTn program was used to query the obtained sequences against the

166    *C. canephora* genome (http://coffee-genome.org/blast). The transcriptome reads (differentially expressed

167    against *H. vastatrix*) were aligned to contig 9 as per Tophat2 (-N 3 --read-gap-length 3 --read-edit-dist 6 -

168    -no-coverage-search --b2-very-sensitive) [33] and to locate the region of the contig containing the genes

169    encoded in response to pathogenicity. The intergenic physical position, distance and orientation were

170    analyzed for the RGAs.

171

## 172    Point mutation analysis

173    The RGAs were analyzed for indels and substitutions using the EMBL MUSCLE multiple

174    sequence alignment tool (http://www.ebi.ac.uk/Tools/msa/muscle/) and MEGA7 [34]. Gene duplication

175    was exclusively analyzed using MEGA 7, while DNA polymorphism and non-synonymous/synonymous

176    substitution rates (ka/ks) were analyzed using DnaSP v5.1 [35].

177

## 178    Functional and phylogenetic analysis

179    Based on the molecular evolution of protein domains, functional diversity between two NBS-

180    LRR RGAs from coffee was analyzed. Homology was compared for the two RGAs and to identify

181    orthologous genes in the genomes of *S. lycopersicum* (https://solgenomics.net/tools/blast/) and *V. vinifera*

182    (https://phytozome.jgi.doe.gov/pz/portal.html). Subsequently, a protein sequence-based comparative

183    phylogenetic tree was constructed for the two genes and their orthologs from the two related genomes

184    using MEGA7 program [34]. The evolutionary history was inferred using the minimum evolutionary

185    method [36].

186

# 187    Results

## 188    Resistance gene screening among the differential coffee clones

189          To investigate the linkage of RGAs to known SH genes, differential coffee clones with different

190     SH genes were subjected to RGA screening using the functional marker, CARF005. Of the 21 differential

191     coffee clones, the marker was detected in eight clones as presented in Table 1 and S2 Fig.  All clones with

192     the $S_H6$ gene had the marker, while those without the gene failed to amplify in the PCR. Thus, gel

193     analysis of the PCR amplicon revealed that this particular RGA marker amplified the $S_H6$ gene locus;

194     however, two exceptions were observed. One of them was that the gene's allele was detected in CIFC

195     128/2-Dilla & Alghe, which is supposed to have just the $S_H1$ gene. In addition, CARF005 was found to be

196     amplified in a differential-host clone CIFC 644/18 H. Kawisari, for which no $S_H$ gene has been reported

197     to date.

198

199

200

201

202

203

204

205

206

207

208

209

210    **Table 1. S$_H$ gene allelic polymorphism detection in 22 differential coffee clones using CARF005 marker.**

| No. | Differential clone* | Susceptible to (*H. vastatrix* physiological race) | S$_H$ gene conferred** | Allelic difference (+/-) |
|---|---|---|---|---|
| 1 | 832/1- Híbrido Timor | - | 6,7,8,9,? | + |
| 2 | HW17/12 | XVI,XXIII | 1,2,4,5 | - |
| 3 | 1343/269- Híbrido Timor | XXII,XXV,XXVI,XXVII,XXVIII,XXIX, XXXI,XXXII,XXXIII,XXXVII,XXXIX,XL | 6 | + |
| 4 | H153/2 | XII, XVI | 1,3,5 | - |
| 5 | H420/10 | XXIX | 5,6,7,9 | + |
| 6 | 110/5-S 4 Agaro | X,XIV,XV, XVI,XXIII,XXIV,XXVI, XXVIII | 4,5 | - |
| 7 | 128/2-Dilla and Alghe | III, X, XII, XVI, XVII, XIX,XXIII, XXVII | 1 | + |
| 8 | 134/4-S12 Kaffa | X, XVI, XIX, XX, XXIII ,XXVII, | 1,4 | - |
| 9 | H419/20 | XXIX, XXXI | 5,6,9 | + |
| 10 | 635/3-S 12 Kaffa | X, XIV,XV,XVI,XIX, XXIII,XXIV,XXVI,XXVII,XXVIII | 1,4,5 | - |
| 11 | 87/1-Geisha | III, X, XII, XVI, XVII, XXIII | 1,5 | - |
| 12 | 1006/10-KP 532 | XII,XVI,XVII, XXIII | 1,2,5 | - |
| 13 | 7963/117-Catimor | XXXIII | 5,7 or 5,7,9 | - |
| 14 | H420/2 | XXIX, XXX | 5,8 | - |
| 15 | 4106 | - | 5,6,7,8,9,? | + |
| 16 | 644/18 H. Kawisari | XIII | ? | + |
| 17 | 832/2- Híbrido Timor | - | 6,7,8,9,? | + |
| 18 | H147/1 | XIV, XVI | 2,3,4,5 | - |
| 19 | 32/1-DK1/6 | I,VIII, XII, XIV, XVI, XVII, XXIII,XXIV, XXV, XXVIII, XXXI | 2,5 | - |

9

| | | | | |
|---|---|---|---|---|
| 20 | H152/3 | XIV,XVI, XXIII, XXIV, XXVII | 2,4,5 | - |
| 21 | 33/1-S.288-23 | VII, VIII, XII, XIV,XVI, | 3,5 | - |
| 22 | Caturra (c) | All | 5 | - |
| 23 | Catuaí 2143-236 (c) | All | 5 | - |
| 24 | Mundo Novo -376/4 (c) | All | 5 | - |

211    *Differential clones were from CIFC (Centro de Investigação das Ferrugens do Cafeeiro, Portugal).

212    **$S_H$1-9 genes as inferred by CIFC.

213    Unknown race (-), coffee genotypes used as negative control (c), presence/absence of allelic differences among $S_H$ genes (+/-) and unknown $S_H$ gene (s) (?).

## Sequence analysis of ORFs

Identification of a BAC clone using CARF005 and the comparative analysis of the RGAs with the other ORFs from *C. canephora* was performed to localize their relative position and determine the putative function. To characterize genetic loci conferring resistance to leaf rust, a BAC clone 78-K-10 (with ~146 kb insert) was identified as shown in S1 b & c Fig. Illumina HiSeq2000/2500 100PE generated 8,711,320 reads. After removing vector sequences and noisy sequences, quality sequences (>20 QC) were assembled into 86 contigs of ≥ 200bp from which the two contigs, contigs 3 (16570 bp) and 9 (8285 bp), were selected (as they had >90% similarity with *Coffea canephora* DNA sequence) and then subjected to downstream processing. The sequences of the two contigs were combined and deposited at NCBI (accession number: KY485942). These contigs shared ≥ 90% identity with *C. canephora* contigs at different chromosomal regions with the highest identity (99% for contig 3 and 97% for contig 9) being on chromosome 0. All the 13 ORFs predicted (eight in contig 3 and five in contig 9) had matched to different species when queried against BLASTp database or to the *C. canephora* genome hub with significant similarities (≤ 1e$^{-05}$ e-value) when BLASTn was used as presented in S3 Table. Among these, five genes (genes 5, 9, 10, 11 and 12) shared significant identities with RGAs from *C. canephora*. These genes are homologous to sequences in the *C. canephora* genome with the highest query coverage being on chromosome 0 as presented in Table 2. Genes 5 (intron-less, 1130 aa) and 11 (with two introns and two exons, 1118 aa) were the largest genes predicted. Both genes were located on the negative reading frames that belong to the CC-NBS-LRR gene family. Mapping to the *C. canephora* genome showed that these genes are separated by 1,634,522 bp, although they are delimited with a shorter length (460 bp) in *C. arabica*. In contrast, the other four RGAs matched and retained their expected positions in the *C. arabica* genome as shown in S4 Fig.

11

243 **Table 2. Size and structure of five resistance gene analogs and their mapping to chromosome 0 of *C. canephora* genome.**

| | Genes[*] | | | | |
|---|---|---|---|---|---|
| | 5 | 9 | 10 | 11 | 12 |
| Contig | 3 | 9 | 9 | 9 | 9 |
| Exon 1 | 3393 | 113 | 121 | 1175 | 345 |
| Intron 1 | - | 554 | 87 | 611 | 1786 |
| Exon 2 | - | 118 | 112 | 2222 | 183 |
| Intron 2 | - | 121 | 711 | 124 | - |
| Exon 3 | - | 69 | 121 | - | - |
| Intron 3 | - | 91 | - | - | - |
| Exon 4 | - | 155 | - | - | - |
| Intron 4 | - | 476 | - | - | - |
| Exon 5 | - | 130 | - | - | - |
| Query coverage (%) | 99.94 | 72.68 | 30.48 | 99.46 | 97.33 |
| Identity (%) | 76.00 | 85.00 | 79.00 | 68.84 | 73.00 |
| E-value | 0.00 | 9,00E-30 | 5,00E-17 | 0.00 | 3,00E-48 |
| Frame | N | N | P | N | P |
| Start hit-End hit | 108638370-108641761 | 106998076-106999730 | 107000654-107000761 | 107000357-107003848 | 107000234-107004551 |
| Protein (aa) | 1130 | 194 | 117 | 1118 | 175 |

244    [*]Exon and intron sizes are in nucleotides.

245    N, negative reading frame and P, positive reading frame.

246    Gene prediction was performed by Augustus command-line version gene prediction [29].

## CARF005 amplicon verification

248    Web-based PCR analysis was conducted to validate the specificity of the CARF005 primer pair and

249    to complement the PCR amplification experiments. *In silico* PCR analysis using contig 9 and gene 11 ORF as

250    the template strands, indicated that the CARF005 marker had a size of 400 bp

251    (http://www.bioinformatics.org/sms2/pcr_products.html). This size of the amplicon was confirmed by PCR

252    using the template DNA from the clone 78-K-10 as outlined in S1 c Fig. Notably, the amplicon spans from

253    position 6867 to 7266 bp in contig 9 (8285 bp) in a negative orientation and from position 2065 to 3115 bp in

254    the ORF of gene 11(3354 bp), respectively.

255

## Gene annotation

257    Gene annotation was performed to identify the putative protein-coding gene. The annotation of 13

258    ORFs showed a range of protein arrays most of which had no role in disease resistance and lacked conserved

259    domains. Among the five RGAs detected in either NCBI BLASTp, or BLASTn against the *C. canephora*

260    genome, genes 9 (unnamed protein product), 10 (putative resistance gene) and 12 (putative resistance gene)

261    had similarity to RGAs as observed by mapping to *C. canephora* genome as presented in S3 Table. Yet, genes

262    5 and 11 encoded the largest resistance proteins (Gene 5: 126.81 kDa and pi: 7.65; gene 11: 126.67 kDa and

263    pi: 8.44) and identified several resistance associated GO terms characterizing their multiple functional

264    domains as shown in Table 3.

265

266    **Table 3. Annotation and functional comparison of gene 5 and 11.**

| Molecular function ontology | | | | |
|---|---|---|---|---|
| **GO ID** | **GO term** | **Reliability (%)** | **Gene 5** | **Gene 11** |
| GO:1901363 | Heterocyclic compound binding | 49 | ✓ | ✓ |
| GO:0000166 | Nucleotide binding | 49 | ✓ | ✓ |
| GO:0005488 | Binding | 49 | ✓ | ✓ |
| GO:1901265 | Nucleoside phosphate binding | 49 | ✓ | ✓ |

| GO:0097159 | Organic cyclic compound binding | 49 | ✓ | ✓ |
|---|---|---|---|---|
| GO:0036094 | Small molecule binding | 49 | ✓ | ✓ |
| GO:0097367 | Carbohydrate derivative binding | 41 | ✓ | ✓ |
| GO:0017076 | Purine nucleotide binding | 41 | ✓ | ✓ |
| GO:0032559 | Adenyl ribonucleotide binding | 41 | ✓ | ✓ |
| GO:0032555 | Purine ribonucleotide binding | 41 | ✓ | ✓ |
| **Biological process ontology** | | | | |
| GO:0006952 | Defense response | 36 | ✓ | ✓ |
| GO:0006950 | Response to stress | 36 | ✓ | ✓ |
| GO:0050896 | Response to stimulus | 36 | ✓ | ✓ |
| GO:0002376 | Immune system process | 16 | ✓ | ✓ |
| GO:0006955 | Immune response | 16 | ✓ | ✓ |
| GO:0045087 | Innate immune response | 16 | ✓ | ✓ |
| GO:0044699 | Single-organism process | 14 | ✓ | ✓ |
| GO:0009987 | Cellular process | 14 | ✓ | ✓ |
| GO:0044763 | Single-organism cellular process | 14 | ✓ | ✓ |
| GO:0033554 | Cellular response to stress | 12 | ✓ | ✓ |
| GO:0016265 | Death | 12 | ✓ | ✓ |
| GO:0051716 | Cellular response to stimulus | 12 | ✓ | ✓ |
| GO:0012501 | Programmed cell death | 12 | ✓ | ✓ |
| GO:0008219 | Cell death | 12 | ✓ | ✓ |
| GO:0034050 | Host programmed cell death induced by symbiont | 12 | ✓ | ✓ |
| GO:0009626 | Plant-type hypersensitive response | 12 | ✓ | ✓ |
| GO:0009814 | Defense response, incompatible interaction | 7 | ✓ | ✓ |
| **Cellular component ontology** | | | | |
| GO:0016020 | Membrane | 33 | ✓ | ✓ |
| GO:0044464 | Cell part | 33 | ✓ | ✓ |
| GO:0005623 | Cell | 33 | ✓ | ✓ |

| GO:0005737 | Cytoplasm | 32 | ✓ | ✓ |
|---|---|---|---|---|
| GO:0044424 | Intracellular part | 32 | ✓ | ✓ |
| GO:0005886 | Plasma membrane | 31 | ✓ | ✓ |
| GO:0071944 | Cell periphery | 31 | ✓ | ✓ |
| GO:0043227 | Membrane-bounded organelle | 24 | ✓ | ✓ |
| GO:0043226 | Organelle | 24 | ✓ | ✓ |
| GO:0005634 | Nucleus | 24 | ✓ | ✓ |

267  Annotation was performed by Predict Protein online server [31] (URL: https://www.predictprotein.org).

268

# Gene characterization

270  To identify the candidate R genes activated against *H. vastatrix* incursion, two contigs (contigs 3 and

271  9) were mapped against the transcriptome of *C. arabica*-*H. vastatrix* interaction [25]. Genes 9, 10, 11 and 12

272  were mapped to transcripts differentially-expressed during incompatible interactions at 12 and 24 hai with *H.*

273  *vastatrix* (race XXXIII) as illustrated in S5 Fig. Contig 3, from which gene 5 was predicted, was also mapped

274  against the same transcriptome resulting in no matching transcripts that were differentially expressed at the

275  two time points (12 and 24 hai) following pathogen inoculation. Contig 9 region, where most R genes are

276  positioned was estimated by mapping against the transcriptome. The result showed approximately 81.58% of

277  the contig encodes transcripts of varying lengths associated with rust resistance, which are activated at 12 and

278  24 hai in response to *H. vastatrix* inoculation. Further analysis of genes 5 and 11 revealed that they belong to

279  the NBS-LRR gene family (the major R genes in plants), suggesting the importance of continuing the *in silico*

280  comparative structural and functional analysis. Intriguingly, both have the Rx-cc-like coiled-coil potato virus

281  x resistance protein domain and four additional multi-domains featuring the entire protein sequence (Fig 1).

282  These genes can be referred as CC-NBS-LRR, as they comprise the N-terminal CC and LRR C-terminal

283  domains flanking the NBS on either side.

284

285  **Fig 1. Comparison of conserved domains and motif architecture in genes 5 and 11.** Note the

286  polymorphism of domains in both genes (encircled by red spheres). Conserved domains were detected using

287  NCBI BLASTp (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins) database.

15

288

289    In addition, annotation of both genes indicates that they encode defense proteins involved in various

290    biological defense as demonstrated in Table 3. Notably, although these genes share 90.24% nucleotide

291    identity, their amino acid sequence identity is only 80.03% (Fig 2). The possibility of substitution mutation

292    events was considered in explaining the diversity. Accordingly, the overall amino acid diversity was

293    attributed to non-synonymous substitution events (non-synonymous/synonymous ratio, ka/ks = 1.5913) in

294    both genes. Further analysis of LRR region showed a higher rate of non-synonymous substitution mutation

295    (ka/ks, non-synonymous/synonymous substitution ratio = 1.9660).

296

297    **Fig 2. Alignment of proteins encoded by genes 5 and 11 and the protein binding regions.** *In silico*

298    prediction of protein binding regions of gene 5 (boxed in red), protein binding regions of gene 11 (boxed in

299    green) and conserved protein binding region (boxed in blue) are shown. Amino acid substitution: unrelated

300    amino acid substitution (space), weakly similar substitution (period), strongly similar substitution (colon) and

301    conserved amino acids (star, unmarked). Note the seven substitution mutations resulting in polymorphism of

302    seven protein binding sites (purple encircled) in either of the sequence. Sequence alignment was carried out

303    using Clustal Omega (http://www.ebi.ac.uk/Tools/msa/clustalo/).

304

# Comparative analysis of structural and functional sites

306    Structural modeling and comparative analyses of the identified genes was performed in order to infer

307    the possible protein functions. Therein, we found that the number of LRR domains in genes 5 and 11 is

308    conserved (13 repeats in both), but arranged differently. We noted the introduction of a coenzyme domain

309    (CoaE, dephospho-CoA kinase) in gene 5, while a LRR variant (LRR_8) seemed to be introduced in gene 11

310    (Fig 1). Despite sharing similar protein multi-domains, two trans-membrane motifs (spanning 5-22 and 102-

311    119 amino acid regions) were detected exclusively in the coiled-coil domain of gene 11 (data not shown). The

312    amino acid sequences of genes 5 and 11 were further analyzed for protein and nucleotide binding site

313    polymorphism. Protein binding sites of the two genes revealed different sensitivity to substitution mutations.

314    Few sites that were specific to each gene were highly affected while most of the binding sites had moderate to

315   no effect as presented in S6 Table. The analysis revealed 17 protein binding sites in gene 5 and 14 sites in

316   gene 11. Similarly, their secondary structures and solvent accessibility properties revealed conserved features

317   (Fig 3aI-IV & cI-IV). Nevertheless, the amino acid residues forming the protein binding sites of the two genes

318   showed high variability in the LRR regions (Fig 2; Fig 3a & c). Although most of the residues are not

319   conserved, the ADP binding site of the NBS domain contained some conserved sites (Fig 3b.II & d.II).

320

321   **Fig 3. *In silico* 3D structure, protein and nucleotide binding site prediction for gene 5 (A and B) and 11**

322   **(C and D).** Protein binding and secondary structure (A and C): Protein binding sites (I), the three types of

323   secondary structures are assumed at different regions (helical: red boxes, strand: blue boxes and loop:

324   intervening white spaces) (II), solvent accessibility (exposed: blue boxes, buried: yellow boxes and

325   intermediate: white spaces) (III) and high protein disorder and flexibility: green boxes (IV). 3D structure and

326   nucleotide binding sites (B and D): 3D structures with Rx-CC-like (blue) to LRR (red to light: orange forming

327   the 'horseshoe' structure) domains are as shown in Figs 3b.I & d.I and the colored residues (NBS) forming

328   the nucleotide (ATP/GTP/ADP/GDP)-binding site (BII and DII). Nucleotide binding site residues with the

329   highest C-score are listed in the right box (conserved residues highlighted in yellow) with the red arrow

330   indicating the sites. I-TASSER modelling C-score [32] was -1.73 and -1.75 (C-score ranging from -5 to 2,

331   where 2 refers to the highest confidence) and 0.29 and 0.17 (C-score ranging from 0-1, where higher score

332   indicates reliable prediction) for nucleotide binding prediction for the two proteins, respectively.

333

334   # Interlocus comparison of $S_H$ genes

335   To investigate the conserved regions in the five RGAs (genes 5, 9, 10, 11 and 12), contigs 3 and 9

336   were queried against three *C. canephora* and 10 *C. arabica* specific contigs assembled from BAC clones

337   spanning $S_H3$ locus from the work of [6]. All the 10 $S_H3$ contigs matched with contig 3 but with varying

338   alignment lengths and identities as presented in S7 Table. Contig GU123898 and HQ696508 (both specific to

339   *C. arabica*) had the highest number of matches to contig 9 (from which four clustered RGAs were predicted)

340   with alignment lengths of 170 nts (77.647% identity and $1.57e^{-31}$ e-value) and 179 nts (76.536% identity and

17

341     1.21e[-26] e-value), respectively. The closest contig (HQ696508) is located on the complementary strand of

342     gene 11 and is 505 bp upstream of the position where CARF005 forward primer annealed to gene 11.

343

## 344 Phylogenetic analysis

345     In an attempt to discern the ancestral relationship of a set of sequences, phylogenetic analysis was

346     performed. Accordingly, two resistance gene families (the NBS-LRR and non-NBS-LRR) were identified,

347     completely sequenced and mapped to chromosome 0 of *C. canephora* genome with a query coverage of

348     99.94% for genes 5 and 11, 72.68% for gene 9, 33.05% for gene 10 and 97.52% for gene 12. The diversity of

349     the NBS-LRR family was detected by analyzing the ka/ks ratios as presented in Table 4. The analysis

350     revealed that the non-synonymous substitution event is common in the CDS, as revealed from all the pairwise

351     analyses. Furthermore, the non-synonymous substitution of CDS is more prominent in the LRR region (in

352     almost all pairwise comparisons) as demonstrated in Table 4.

353

354 **Table 4. Pair-wise synonymous and non-synonymous nucleotide substitution analysis among the six**
355 **resistance gene analogs (gene 5 and 11 and their respective two top hits as mined by BLASTn in NCBI).**

| | | Entire protein | | | LRR region | | |
|---|---|---|---|---|---|---|---|
| Seq. 1 | Seq. 2 | Ks | Ka | ka/ks | Ks | Ka | ka/ks |
| gene5_hit1 | gene11_hit1 | 0.0786 | 0.1302 | 1.6565 | 0.0702 | 0.1383 | 1.9701 |
| gene5_hit1 | gene11_hit2 | 0.0899 | 0.1614 | 1.7953 | 0.0536 | 0.1408 | 2.6269 |
| gene5_hit1 | gene11 | 0.0723 | 0.1177 | 1.6279 | 0.0622 | 0.1233 | 1.9823 |
| gene5_hit1 | gene5_hit2 | 0.0635 | 0.0999 | 1.5732 | 0.0583 | 0.1029 | 1.7650 |
| gene5_hit1 | gene5 | 0.0009 | 0.0039 | 4.3333 | 0.0015 | 0.0045 | 3.0000 |
| gene11_hit1 | gene11_hit2 | 0.1092 | 0.1854 | 1.6978 | 0.0756 | 0.1602 | 2.1190 |
| gene11_hit1 | gene11 | 0.0061 | 0.0164 | 2.6885 | 0.0095 | 0.0170 | 1.7895 |
| gene11_hit1 | gene5_hit2 | 0.0761 | 0.1309 | 1.7201 | 0.0736 | 0.1369 | 1.8601 |
| gene11_hit1 | gene5 | 0.0786 | 0.1288 | 1.6387 | 0.0701 | 0.1383 | 1.9729 |
| gene11_hit2 | gene11 | 0.1074 | 0.1742 | 1.6220 | 0.0686 | 0.1445 | 2.1064 |
| gene11_hit2 | gene5_hit2 | 0.0846 | 0.5829 | 6.8901 | 0.0607 | 0.1440 | 2.3723 |
| gene11_hit2 | gene5 | 0.0902 | 0.1620 | 1.7960 | 0.0540 | 0.1430 | 2.6481 |
| gene11 | gene5_hit2 | 0.0704 | 0.1155 | 1.6406 | 0.0616 | 0.1199 | 1.9464 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| gene11 | gene5 | 0.0723 | 0.1164 | 1.6100 | 0.0622 | 0.1234 | 1.9839 |
| gene5_hit2 | gene5 | 0.0635 | 0.0997 | 1.5701 | 0.0583 | 0.1052 | 1.8045 |

356 Seq., sequence, non-synonymous/synonymous substitution rate was computed by DnaSP v5.1 [35].

357

358 Moreover, phylogenetic analysis showed that the tomato gene 5 was closely related to genes 5 and

359 11 of coffee than the gene 11 of both tomato and grape (Fig 4). Within coffee itself, a significant diversity

360 between genes 5 and 11 was detected by the MEGA 7 bootstrap method of the phylogenetic test.

361

362 **Fig 4. Phylogenetic history of genes 5 and 11 in three related genomes.** The evolutionary history was

363 inferred using the Minimum Evolution method [36]. The optimal tree with the sum of branch length =

364 2.98805978 is shown. The percentage of replicate trees with the associated taxa clustered together in the

365 bootstrap test (500 replicates) are shown next to the branches [67]. The tree is drawn to scale, with branch

366 lengths in the same units as those of the evolutionary distances used to infer the phylogenetic tree. The

367 evolutionary distances were computed using the Poisson correction method [68] and are in the units of the

368 number of amino acid substitutions per site. The ME tree was searched using the Close-Neighbor-Interchange

369 (CNI) algorithm [69] at a search level of 1. The Neighbor-joining algorithm [70] was used to generate the

370 initial tree. The analysis involved 6 amino acid sequences. All positions containing gaps and missing data

371 were eliminated. A total of 554 positions were there in the final dataset. Evolutionary analyses were

372 conducted in MEGA7 [34]. Subject IDs are indicated in parenthesis for the corresponding two homologous

373 sequences mined by BLASTx against tomato (Sol Genomics Network) and grape (Phytozome) genome

374 databases.

375

# Discussion

377 The majority of NBS-LRR encoding genes are known to be clustered but unevenly distributed in

378 plant genomes [10,37–39]. The NBS-ARC domain is known to be involved in directly blocking the biotrophic

379 pathogens by activating the hypersensitive response (HR) [40]. HR starts with programmed cell death of

380 affected and surrounding cells and ends with the activation of systemic acquired resistance (SAR), in which

381     the defense is induced in distal non-infected cells of the host under attack [41,42]. By recognizing the

382     corresponding virulence (vr) factors or their effects, NBS-LRR proteins are sufficient to induce HR

383     [8,9,42,43]. In the present study, a cluster of two different classes of RGAs resistant to coffee rust, the NBS-

384     LRRs linked to non-NBS-LRR genes were reported. The two NBS-LRR genes (genes 5 and 11) are the

385     largest non-TNL genes sequenced in Arabica coffee and most other plants investigated to date [6,44–46].

386     We identified 13 genes: genes 1, 2, 3, 4, 5, 6, 7 and 8 (gene 5 is a RGA) from contig 3 and genes 9,

387     10, 11, 12, and 13 (only 13 is not a RGA) from contig 9. We also report, a completely sequenced and

388     characterized novel RGA (gene 11), from Híbrido de Timor CIFC 832/2, probably a major component of the

389     $S_H$ gene. This Híbrido de Timor (HDT) (*C. arabica* x *C. canephora*) is immune to all known virulence factors

390     of *H. vastatrix* physiological races and therefore, is an extremely important source of resistance [47].

391     Additionally, the mapping against the transcriptome of *C. arabica-H. vastatrix* interaction suggests that the

392     three other clustered RGAs (genes 9, 10 and 12) have differential-expression during the incompatible

393     interaction. Mapping study has also revealed the presence of reads exclusively mapped to transcripts of

394     pathogen-infected plants at 12 and 24 hai.

395     The known rust resistance genes, $S_H3$ in *C. liberica* [48], $S_H6$, 7, 8 and 9 in *C. canephora* [49] and

396     $S_H1$, 2, 4 and 5 genes in *C. arabica* are dominantly-inherited genes [50]. One of the fundamental questions to

397     be clarified is how different are these 9 $S_H$ genes that belong to different coffee species. The comparative

398     analysis of contigs from the $S_H3$ locus of *C. arabica* and *C. canephora* [6] revealed different levels of

399     conservation of motifs in the two contigs examined: contigs 3 and 9. The results indicate that the RGAs may

400     share large conserved regions, but few highly polymorphic regions encoding specific protein motifs necessary

401     for critical roles. This characteristic conservation of domains was once more confirmed based on comparative

402     analysis of the cloned gene (gene 11) and using differential coffee clones for $S_H$ gene identification. PCR

403     amplification of gene 11 also indicated the existence of allelic difference/polymorphism among the $S_H$ gene

404     loci and considerable sequences of conserved domain (on which CARF005 primer was designed) with $S_H6$

405     and possibly with $S_H1$. As PCR amplification using CARF005 primer (constituting gene 11) was detected in

406     all the differential clones with $S_H6$ and 832/1-HT and 832/2-HT containing $S_H6$, 7, 8, 9 and $S_H$?. In addition,

407     we report a conserved sequence of gene 11 11 in CIFC 128/2-Dilla & Alghe, previously considered to contain

408     only the $S_H1$ gene, and in CIFC 644/18 H. Kawisari with an uncharacterized $S_H$-gene. Overall, we propose

409 the following hypothesis for extensive and rigorous biological investigation: the identified gene (gene 11)

410 could be one of the unidentified and a not yet supplanted (at least in Brazil) $S_H$ gene in HDT consisting of a

411 conserved domain (CARF005) shared with the $S_H6$ and $S_H1$ genes.

412     Mapping of the RGAs to *C. canephora*, the result from differential clone screening and annotation

413 altogether confirmed that gene 11 locus is descended from *C. canephora,* hence is a sibling of $S_H6$-9 [49].

414 Besides, mapping of RGAs to *C. canephora* was complemented by the differential clone screening for SH

415 genes, which affirms the existence of strong linkage of SH gene locus and the RGAs as their conserved

416 sequences (CARF005) were detected in eight of the differential clones. The disparity of the position of gene 5

417 in relation to gene 9 (the fact that all the predicted genes are from an insert size of ~146 kbp) could be

418 attributed to the linked LTR-retrotransposon (GROUP_78_RLC4) as demonstrated in  S4 Fig and the

419 transposase gene (gene 1). Transposons could have interrupted and separated the two genes apart by 1.6 Mb,

420 since *C. arabica* is known to have diverged from *C. canephora*. Multiple transposable elements linked to

421 NBS-LRR regions were reported in other plants [45,51]. Transposition of genes and gene fragments are some

422 of the mechanisms that generate variability and positional changes among the NBS-LRR genes in different

423 plants [45,51–54]

424     Rx-CC, PLN and NB-ARC domains are conserved in the NBS-LRR genes across diverse plant

425 species [44,55,56]. The potato virus x resistance (Rx) protein-like N-terminal coiled-coil domain mediates

426 intramolecular interaction with NB-ARC and intermolecular interactions through RanGAP2 (Ran-GTPase-

427 activating protein-2) in potato [43,57]. Rx-CC, RanGAP2 interaction site and NB-ARC were detected in

428 genes 5 and 11, suggesting similarity in their defense role in coffee. However, unlike the Rx-CC domain with

429 four helical structures, five helical structures are conserved in genes 5 and 11, indicating polymorphic

430 differences between the species. The PLN00113 domain in gene 5 and PLN03210 in gene 11, span the LRR

431 region and were initially reported in *A. thaliana* [44]. The distinct position of these domains in genes 5 and 11

432 indicates high variability in the LRR region in both genes. Functional motif prediction indicated that the

433 PLN03210 (LRR domain) is likely engaged in direct effector interaction while the corresponding PLN00113

434 of gene 5 is engaged in LRR-reception and downstream kinase-mediated signaling. These observations are in

435 accordance with the functional and structural analysis data of LRR proteins in *A. thaliana* [44,58–61]. Based

436 of their annotations, the two genes (genes 5 and 11) products are intracellular resistance proteins that directly

437  or indirectly recognize pathogen effector proteins and subsequently trigger a response that may be as severe

438  as localized cell death [42].

439      Different selection pressures shape the evolution of domains in the NBS-LRR encoding genes. The

440  NBS domain was assumed to be under the purifying selection (a negative selection in which variation is

441  minimized by stabilizing selection) than by the diversifying selection, which acts on the LRR domain [9,62].

442  In contrast, the diversifying selection (positive selection) act on all the domains of genes 5 and 11 (ka/ks >1).

443  This result is contrary to the general assumption that diversifying selection is diluted when the overall non-

444  synonymous substitution is considered [6], indicating an intense diversifying selection action on both genes.

445  Further investigation of four more orthologous genes also resulted in similar findings, indicating that the

446  NBS-LRR genes are highly variable due to substitution mutations. As the LRR domains are involved in direct

447  ligand binding, their variability due to non-synonymous substitution is higher than that seen in other domains.

448  This results in the formation of a super-polymorphic region to cope with the continuously evolving pathogen

449  effectors. Similar findings (from different plants, including coffee) on diversifying selection have been

450  reported [6,9,11,38,45,63,64]. Diversifying selection by non-synonymous substitution was detected in non-

451  NBS-LRR genes (genes 10 and 12) (data not shown), reiterating the importance of substitution mutation in

452  such clustered R genes. Synergistic activation of the two groups (NBS-LRR and non-NBS-LRR) may

453  enhance the resistance durability; and so their expression pattern merits further investigation.

454      Based on the phylogenetic tree of orthologous genes originated from related genomes, the six genes

455  could be divided into two groups. Gene 5 from tomato is closely related to genes 5 and 11 from coffee,

456  making the first group, whereas genes 5 and 11 from grape happens to be the second highly diversified group.

457  Intraspecies diversity of non-TIR-NBS-LRR due to substitution and genetic recombination exist in grape [65]

458  and tomato [66] while gene duplication and conversion events were observed in coffee [6]. In general, the

459  phylogenetic tree revealed that genes 5 and 11 may have recently diverged in coffee, while the divergence

460  observed in the other species may have been earlier events.

461

462  # Conclusion

463        The two groups of RGAs, NBS-LRR and non-NBS-LRR, are clustered in a single locus from which

464    multiple variants of resistance genes are expressed to confer specific resistance functions. The four cloned,

465    sequenced and characterized RGAs span a rust resistance gene locus descended from *C. canephora*. The two

466    CC-NBS-LRR protein encoding genes are under strong diversifying selection impacting all component

467    domains. A more intense diversification of LRR region indicates that the variability in the effector binding

468    site is the cause of divergence in resistance specificity. Although conserved sequences were detected for the

469    $S_H6$ gene across the various differential coffee clones, it could be inferred that the $S_H$ gene loci have a

470    characteristic polymorphism conferring different resistance phenotypes against coffee leaf rust. This is the

471    first report unveiling new insights into the molecular nature of $S_H$ genes. The CC-NBS-LRR gene thus

472    characterized is the largest and most complete sequence ever reported in Arabica coffee. The work

473    demonstrated a cluster of resistance genes spanning the R gene locus that could serve as functional markers

474    for subsequent functional analysis. These findings could also serve as a benchmark for validation of

475    expression patterns in response to pathogenicity and gene segregation along generations. Such studies can be

476    applied in molecular breeding as it has the potential to replace arbitrary DNA-based marker-assisted breeding

477    at least for two reasons. First, there is no loss due to segregation, which is the case even for finely saturated

478    markers. Second, four of the RGAs (genes 9, 10, 11 and 12) are stacked in a locus, from which different

479    primers can be designed to screen genotypes to verify co-segregation analysis.

480

# 481 Supporting information

482    **S1 Fig. Work flow in BAC clone screening.** Clone pooling and subsequent group decomposition to isolate a

483    single clone with CARF005 insert (A), DNA of isolated clone 78-K-10 (B) and CARF005 PCR amplicon (C)

484    as revealed by 1% UltraPure™ agarose gel electrophoresis. M is 100 bp DNA size marker.

485

486    **S2 Fig. The 21 differential coffee clones screened for CARF005 marker (listed in order as in Table 1).**

487    Clones with CARF005 were 1 (832/1-HT), 3 (1343/269-HT), 5 (H420/10), 7 (128/2-Dilla and Alghe), 9

488    (H419/20), 15 (4106), 16 (644/18 H. Kawisan, a new report) and 17 (832/2-HT). M: DNA weight marker

489    ladder (the lightest band being 100 bp). No gel cropping was performed to any of the lanes displayed in the

490    gel above.

491

492    **S3 Fig. Mapping of RGAs clustered on chromosome 0 of *C. canephora*.** Putative mRNA transcription

493    orientations are shown by black arrows. The relatively larger size of gene 12 is due to the largest size of its

494    intron 1 (Table 2). Green boxes are used to mark query positions relative to subject genes (other gene

495    products, all in blue boxes). Note that gene 10 and 12 are in positive orientation (Table 2) with no matching

496    transcript here, hence probably originated from different parent or attributed to mutation events in *C. arabica*.

497    Mapping was carried out by CDS (coding sequence) BLASTn followed by track assembly on *C. canephora*

498    genome hub server [5].

499

500    **S4 Fig. Mapping of contig 9 to transcriptome of differentially expressed genes during *C. arabica-H.***

501    ***vastatrix* (race XXXIII) incompatible interaction to show the region of active gene (gene 9, 10, 11 and**

502    **12) expression.** Note the three expression profiles (three rows) corresponding to control (uninoculated at 0

503    hour, top row), 12 (middle row) and 24 hai (bottom row) of transcriptome reads mapped against contig 9 of

504    resistant coffee clone (CIFC HDT 832/2). Grey shades indicate matching transcriptome reads while

505    nucleotide substitutions (mismatches) were shown by colored strips (yellow: G, green: A, red: R and blue: C).

506    Large red shades indicate deletions. Contig mapping was performed by Tophat 2[28] setting alignment

507    parameter as '-N 3 --read-gap-length 3 --read-edit-dist 6 --no-coverage-search --b2-very-sensitive' to locate

508    the region of the contig encoding genes against the pathogen and visualized with Integrative Genomics

509    Viewer (IGV) v. 2.3[33].

510

511    **S1 Table. Top hits for the 13 ORFs as found in NCBI by BLASTp or at *C. canephora* genome by**

512    **BLASTn.** *Homologous sequences for which no ID/Accession number has been assigned are indicated in

513    hyphen.          BLASTp          was          performed          by          NCBI          online          server

514    (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins).

515

516    **S2 Table. Mutation (substitution) effect on protein binding regions of gene 5 and 11 indicated by amino**

517    **acid sequence in respective genes.** *Hyphen indicates range of amino acids constituting the binding site.

24

518    Yellow highlighted residues are conserved residues in both genes while purple highlighted residues are

519    specific protein binding sites in respective gene. Substitution mutation effect analysis was performed by The

520    Predict Protein Server[31].

521

522    **S3 Table. Output of the two contigs BLASTed against $S_H3$ locus contigs specific to C. arabica and C.**

523    **canephora*.** *Ten contigs specific to *C. arabica* and three contigs specific to *C. canephora*, all assembled

524    from BAC clones with $S_H3$ locus were taken from the work of [6].

525

# Acknowledgements

527    We thank Drs Jorge L.B. Pacheco, Abraham Abera, Bayissa Chala and Mohammed Naimuddin for

528    their valuable suggestions and edition of the manuscript. We are also grateful to the Agronomic Institute of

529    Paraná, Londrina-Brazil, for providing the CIFC 832/2 BAC library.

530

# Author contributions

532    **Conceptualization:** GDB ETC

533    **Formal analysis:** GDB SSF

534    **Investigation:** GDB

535    **Methodology:** GDB ETC SSF

536    **Project administration:** ETC LZ

537    **Resources:** ETC LZ

538    **Software:** GDB SSF

539    **Supervision:** ETC LZ

540    **Validation:** GDB ETC SSF LZ

541    **Writing-Original draft:** GDB

542    **Writing-Review & editing:** GDB ETC SSF LZ

543

# References

544

545  1.    Mussatto SI, Machado EMS, Martins S, Teixeira JA. Production, Composition, and Application of

546       Coffee and Its Industrial Residues. Food Bioprocess Technol. 2011;4: 661–672. doi:10.1007/s11947-

547       011-0565-z

548  2.    Zambolim L. Current status and management of coffee leaf rust in Brazil. Trop Plant Pathol. 2016;41:

549       1–8. doi:10.1007/s40858-016-0065-9

550  3.    Gichuru EK, Ithiru JM, Silva MC, Pereira AP, Varzea VMP. Additional physiological races of coffee

551       leaf rust (Hemileia vastatrix) identified in Kenya. Trop Plant Pathol. 2012;37: 424–427.

552       doi:10.1590/S1982-56762012000600008

553  4.    Cristancho MA, Botero-Rozo DO, Giraldo W, Tabima J, Riaño-Pachón DM, Escobar C, et al.

554       Annotation of a hybrid partial genome of the coffee rust (Hemileia vastatrix) contributes to the gene

555       repertoire catalog of the Pucciniales. Front Plant Sci. 2014;5. doi:10.3389/fpls.2014.00594

556  5.    Denoeud F, Carretero-Paulet L, Dereeper A, Droc G, Guyot R, Pietrella M, et al. The coffee genome

557       provides insight into the convergent evolution of caffeine biosynthesis. Science (80- ). 2014;345:

558       1181–1184. doi:10.1126/science.1255274

559  6.    Ribas AF, Cenci A, Combes M-C, Etienne H, Lashermes P. Organization and molecular evolution of

560       a disease-resistance gene cluster in coffee trees. BMC Genomics. 2011;12: 240. doi:10.1186/1471-

561       2164-12-240

562  7.    Cenci A, Combes M-C, Lashermes P. Comparative sequence analyses indicate that Coffea (Asterids)

563       and Vitis (Rosids) derive from the same paleo-hexaploid ancestral genome. Mol Genet Genomics.

564       2010;283: 493–501. doi:10.1007/s00438-010-0534-7

565  8.    Jones JDG, Dangl JL. The plant immune system. Nature. 2006;444: 323–329.

566       doi:10.1038/nature05286

567  9.    McHale L, Tan X, Koehl P, Michelmore RW. Plant NBS-LRR proteins: adaptable guards. Genome

568       Biol. 2006;7: 212. doi:10.1186/gb-2006-7-4-212

569    10.    Hulbert SH, Webb CA, Smith SM, Sun Q. Resistance gene complexes: evolution and utilization.

570           Annu Rev Phytopathol. 2001;39: 285–312. doi:10.1146/annurev.phyto.39.1.285

571    11.    DeYoung BJ, Innes RW. Plant NBS-LRR proteins in pathogen sensing and host defense. Nat

572           Immunol. 2006;7: 1243–1249. doi:10.1038/ni1410

573    12.    Liu J, Liu X, Dai L, Wang G. Recent Progress in Elucidating the Structure, Function and Evolution of

574           Disease Resistance Genes in Plants. J Genet Genomics. 2007;34: 765–776. doi:10.1016/S1673-

575           8527(07)60087-3

576    13.    Leipe DD, Koonin E V., Aravind L. STAND, a Class of P-Loop NTPases Including Animal and Plant

577           Regulators of Programmed Cell Death: Multiple, Complex Domain Architectures, Unusual Phyletic

578           Patterns, and Evolution by Horizontal Gene Transfer. J Mol Biol. 2004;343: 1–28.

579           doi:10.1016/j.jmb.2004.08.023

580    14.    Van der Hoorn RAL. Identification of Distinct Specificity Determinants in Resistance Protein Cf-4

581           Allows Construction of a Cf-9 Mutant That Confers Recognition of Avirulence Protein AVR4.

582           PLANT CELL ONLINE. 2001;13: 273–285. doi:10.1105/tpc.13.2.273

583    15.    Kushalappa AC, Yogendra KN, Karre S. Plant Innate Immune Response: Qualitative and Quantitative

584           Resistance. CRC Crit Rev Plant Sci. 2016;35: 38–55. doi:10.1080/07352689.2016.1148980

585    16.    Ellis J, Dodds P, Pryor T. Structure, function and evolution of plant disease resistance genes. Curr

586           Opin Plant Biol. 2000;3: 278–284. doi:10.1016/S1369-5266(00)00080-7

587    17.    Combes M-C, Lashermes P, Chalhoub B, Noir S, Patheyron S. Construction and characterisation of a

588           BAC library for genome analysis of the allotetraploid coffee species ( Coffea arabica L.). TAG Theor

589           Appl Genet. 2004;109: 225–230. doi:10.1007/s00122-004-1604-1

590    18.    Cação SMB, Silva N V., Domingues DS, Vieira LGE, Diniz LEC, Vinecky F, et al. Construction and

591           characterization of a BAC library from the Coffea arabica genotype Timor Hybrid CIFC 832/2.

592           Genetica. 2013;141: 217–226. doi:10.1007/s10709-013-9720-y

593    19.    Poczai P, Varga I, Laos M, Cseh A, Bell N, Valkonen JP, et al. Advances in plant gene-targeted and

594            functional markers: a review. Plant Methods. 2013;9: 6. doi:10.1186/1746-4811-9-6

595    20.    Noir S, Combes MC, Anthony F, Lashermes P. Origin, diversity and evolution of NBS-type disease-

596            resistance gene homologues in coffee trees (Coffea L.). Mol Genet Genomics. 2001;265: 654–62.

597            Available: http://www.ncbi.nlm.nih.gov/pubmed/11459185

598    21.    Rodrigues CJ, Bettencourt AJ, Rijo L. Races of the Pathogen and Resistance to Coffee Rust. Annu

599            Rev Phytopathol. 1975;13: 49–70. doi:10.1146/annurev.py.13.090175.000405

600    22.    Alvarenga MS, Caixeta TE, Hufnagel B, Thiebaut F, Maciel-Zambolim E, Zambolim L, et al. In silico

601            identification of coffee genome expressed sequences potentially associated with resistance to diseases.

602            Genet Mol Biol. 2010;33: 795–806. Available: http://www.scielo.br/pdf/gmb/v33n4/31.pdf

603    23.    Van der Biezen EA, Jones JD. Plant disease-resistance proteins and the gene-for-gene concept. Trends

604            Biochem Sci. 1998;23: 454–6. Available: http://www.ncbi.nlm.nih.gov/pubmed/9868361

605    24.    Diniz LEC, Sakiyama NS, Lashermes P, Caixeta ET, Oliveira ACB, Zambolim EM, et al. Analysis of

606            AFLP markers associated to the Mex-1 resistance locus in Icatu progenies. Crop Breed Appl

607            Biotechnol. 2005;5: 387–393. doi:10.12702/1984-7033.v05n04a03

608    25.    Florez JC, Mofatto LS, do Livramento Freitas-Lopes R, Ferreira SS, Zambolim EM, Carazzolle MF,

609            et al. High throughput transcriptome analysis of coffee reveals prehaustorial resistance in response to

610            Hemileia vastatrix infection. Plant Mol Biol. 2017; doi:10.1007/s11103-017-0676-7

611    26.    Alvarenga M. MARCADORES MOLECULARES DERIVADOS DE SEQÜÊNCIAS EXPRESSAS

612            DO GENOMA DO CAFÉ POTENCIALMENTE ENVOLVIDAS NA RESISTÊNCIA À

613            FERRUGEM. Universidade Federal de Vicosa. 2007.

614    27.    Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A New

615            Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. J Comput Biol.

616            2012;19: 455–477. doi:10.1089/cmb.2012.0021

617    28.    Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of

618            transcriptomes in the presence of insertions, deletions and gene fusions. Genome Biol. 2013;14: R36.

619  doi:10.1186/gb-2013-14-4-r36

620  29.  Stanke M, Steinkamp R, Waack S, Morgenstern B. AUGUSTUS: a web server for gene finding in

621  eukaryotes. Nucleic Acids Res. 2004;32: W309–W312. doi:10.1093/nar/gkh379

622  30.  Lin C, Mueller LA, Carthy JM, Crouzillat D, Pétiard V, Tanksley SD. Coffee and tomato share

623  common gene repertoires as revealed by deep sequencing of seed and cherry transcripts. Theor Appl

624  Genet. 2005;112: 114–130. doi:10.1007/s00122-005-0112-2

625  31.  Yachdav G, Kloppmann E, Kajan L, Hecht M, Goldberg T, Hamp T, et al. PredictProtein--an open

626  resource for online prediction of protein structural and functional features. Nucleic Acids Res.

627  2014;42: W337–W343. doi:10.1093/nar/gku366

628  32.  Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y. The I-TASSER Suite: protein structure and

629  function prediction. Nat Methods. 2014;12: 7–8. doi:10.1038/nmeth.3213

630  33.  Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al. Integrative

631  genomics viewer. Nat Biotechnol. 2011;29: 24–26. doi:10.1038/nbt.1754

632  34.  Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis Version 7.0 for

633  Bigger Datasets. Mol Biol Evol. 2016;33: 1870–1874. doi:10.1093/molbev/msw054

634  35.  Rozas J. DNA sequence polymorphism analysis using DnaSP. Methods in molecular biology (Clifton,

635  N.J.). 2009. pp. 337–350. doi:10.1007/978-1-59745-251-9_17

636  36.  Rzhetsky  a, Nei M. A Simple Method for Estimating and Testing Minimum-Evolution Trees. Mol

637  Biol Evol. 1992;9: 945–967.

638  37.  Meyers BC, Kaushik S, Nandety RS. Evolving disease resistance genes. Curr Opin Plant Biol.

639  2005;8: 129–134. doi:10.1016/j.pbi.2005.01.002

640  38.  Hammond-Kosack KE, Kanyuka K. Resistance Genes ( R Genes) in Plants. eLS. Chichester, UK:

641  John Wiley & Sons, Ltd; 2007. doi:10.1002/9780470015902.a0020119

642  39.  Goyal N, Bhatia G, Sharma S, Garewal N, Upadhyay A, Upadhyay SK, et al. Genome-wide

643  characterization revealed role of NBS-LRR genes during powdery mildew infection in Vitis vinifera.

644  Genomics. 2019; doi:10.1016/j.ygeno.2019.02.011

645 40. Mur LAJ, Kenton P, Lloyd AJ, Ougham H, Prats E. The hypersensitive response; the centenary is

646  upon us but how much do we know? J Exp Bot. 2008;59: 501–520. doi:10.1093/jxb/erm239

647 41. Sanabria N, Goring D, Nürnberger T, Dubery I. Self/nonself perception and recognition mechanisms

648  in plants: a comparison of self-incompatibility and innate immunity. New Phytol. 2008;178: 503–514.

649  doi:10.1111/j.1469-8137.2008.02403.x

650 42. Qi D, Innes RW. Recent Advances in Plant NLR Structure, Function, Localization, and Signaling.

651  Front Immunol. 2013;4. doi:10.3389/fimmu.2013.00348

652 43. Rairdan GJ, Collier SM, Sacco MA, Baldwin TT, Boettrich T, Moffett P. The Coiled-Coil and

653  Nucleotide Binding Domains of the Potato Rx Disease Resistance Protein Function in Pathogen

654  Recognition and Signaling. PLANT CELL ONLINE. 2008;20: 739–751. doi:10.1105/tpc.107.056036

655 44. Kim SH, Kwon SI, Saha D, Anyanwu NC, Gassmann W. Resistance to the Pseudomonas syringae

656  Effector HopA1 Is Governed by the TIR-NBS-LRR Protein RPS6 and Is Enhanced by Mutations in

657  SRFR1. PLANT Physiol. 2009;150: 1723–1732. doi:10.1104/pp.109.139238

658 45. Ratnaparkhe MB, Wang X, Li J, Compton RO, Rainville LK, Lemke C, et al. Comparative analysis of

659  peanut NBS-LRR gene clusters suggests evolutionary innovation among duplicated domains and

660  erosion of gene microsynteny. New Phytol. 2011;192: 164–178. doi:10.1111/j.1469-

661  8137.2011.03800.x

662 46. Djebbi S, Bouktila D, Makni H, Makni M, Mezghani-Khemakhem M. Identification and

663  characterization of novel NBS-LRR resistance gene analogues from the pea. Genet Mol Res. 2015;14:

664  6419–6428. doi:10.4238/2015.June.11.18

665 47. Bettencourt J. Considerações gerais sobre o "Híbrido de Timor" [Internet]. 1st ed. Journal of

666  Chemical Information and Modeling. Sao Paulo: Instituto Agronômico; 1973.

667  doi:10.1017/CBO9781107415324.004

668     48.     Noronha-Wagner and Bettencourt A. Genetic study of the resistance of Coffea sp to leaf rust 1.

669            Identification and behavior of four factors conditioning disease reaction in Coffea arabica to twelve

670            physiologic races of Hemileia vastatrix. Can J Bot. 1967;45: 2021–2031.

671     49.     Bettencourt, Rodrigues. Principles and practice of coffee breeding for resistance to rust and other

672            diseases. Elsevier Appl Sci. 1988;3: 199–234.

673     50.     Bettencourt AJ, Coronha-Wagner. Genetic factors conditioning resistance of Coffea arabica L. to

674            Hemileia vastatrix Berk et Br. Agron Lusit. 1971;31: 285–292.

675     51.     Kang Y, Kim K, Shim S, Yoon M, Sun S, Kim M, et al. Genome-wide mapping of NBS-LRR genes

676            and their association with disease resistance in soybean. BMC Plant Biol. 2012;12: 139.

677            doi:10.1186/1471-2229-12-139

678     52.     González VM, Aventín N, Centeno E, Puigdomènech P. Interspecific and intraspecific gene

679            variability in a 1-Mb region containing the highest density of NBS-LRR genes found in the melon

680            genome. BMC Genomics. 2014;15: 1131. doi:10.1186/1471-2164-15-1131

681     53.     Sanseverino W, Hénaff E, Vives C, Pinosio S, Burgos-Paz W, Morgante M, et al. Transposon

682            Insertions, Structural Variations, and SNPs Contribute to the Evolution of the Melon Genome. Mol

683            Biol Evol. 2015;32: 2760–2774. doi:10.1093/molbev/msv152

684     54.     Panchy N, Lehti-Shiu MD, Shiu S-H. Evolution of gene duplication in plants. Plant Physiol. 2016;

685            pp.00523.2016. doi:10.1104/pp.16.00523

686     55.     van der Biezen EA, Jones JD. The NB-ARC domain: a novel signalling motif shared by plant

687            resistance gene products and regulators of cell death in animals. Curr Biol. 1998;8: R226-7.

688            doi:10.1016/S0960-9822(98)70145-9 showArticle Info

689     56.     Wang G-F, Ji J, EI-Kasmi F, Dangl JL, Johal G, Balint-Kurti PJ. Molecular and Functional Analyses

690            of a Maize Autoactive NB-LRR Protein Identify Precise Structural Requirements for Activity.

691            Mackey D, editor. PLOS Pathog. 2015;11: e1004674. doi:10.1371/journal.ppat.1004674

692     57.     Hao W, Collier SM, Moffett P, Chai J. Structural Basis for the Interaction between the Potato Virus X

693        Resistance Protein (Rx) and Its Cofactor Ran GTPase-activating Protein 2 (RanGAP2). J Biol Chem.

694        2013;288: 35868–35876. doi:10.1074/jbc.M113.517417

695   58.   Lahaye T. The Arabidopsis RRS1-R disease resistance gene--uncovering the plant's nucleus as the

696        new battlefield of plant defense? Trends Plant Sci. 2002;7: 425–7. Available:

697        http://www.ncbi.nlm.nih.gov/pubmed/12399170

698   59.   Kierszniowska S, Seiwert B, Schulze WX. Definition of Arabidopsis Sterol-rich Membrane

699        Microdomains by Differential Treatment with Methyl- -cyclodextrin and Quantitative Proteomics.

700        Mol Cell Proteomics. 2009;8: 612–623. doi:10.1074/mcp.M800346-MCP200

701   60.   Gou X, He K, Yang H, Yuan T, Lin H, Clouse SD, et al. Genome-wide cloning and sequence analysis

702        of leucine-rich repeat receptor-like protein kinase genes in Arabidopsis thaliana. BMC Genomics.

703        2010;11: 19. doi:10.1186/1471-2164-11-19

704   61.   Xu Y, Liu F, Zhu S, Li X. The Maize NBS-LRR Gene ZmNBS25 Enhances Disease Resistance in

705        Rice and Arabidopsis. Front Plant Sci. 2018;9. doi:10.3389/fpls.2018.01033

706   62.   Michelmore RW, Meyers BC. Clusters of resistance genes in plants evolve by divergent selection and

707        a birth-and-death process. Genome Res. 1998;8: 1113–30. Available:

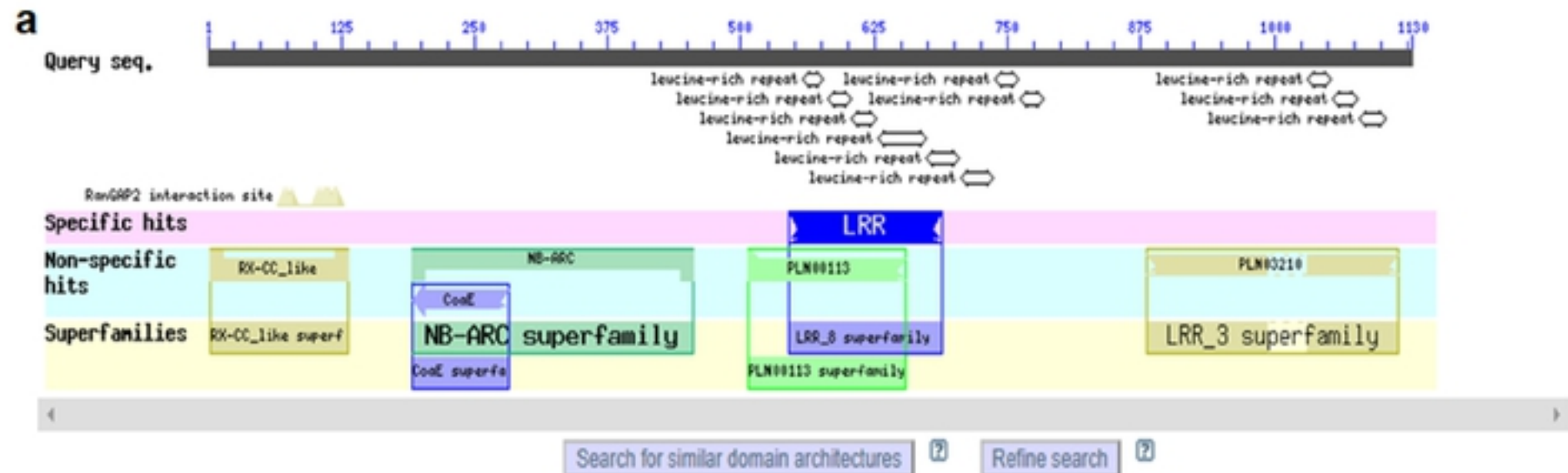708        http://www.ncbi.nlm.nih.gov/pubmed/9847076

709   63.   Padmanabhan M, Cournoyer P, Dinesh-Kumar SP. The leucine-rich repeat domain in plant innate

710        immunity: a wealth of possibilities. Cell Microbiol. 2009;11: 191–198. doi:10.1111/j.1462-

711        5822.2008.01260.x

712   64.   Zhao Y, Huang J, Wang Z, Jing S, Wang Y, Ouyang Y, et al. Allelic diversity in an NLR gene BPH9

713        enables rice to combat planthopper variation. Proc Natl Acad Sci. 2016;113: 12850–12855.

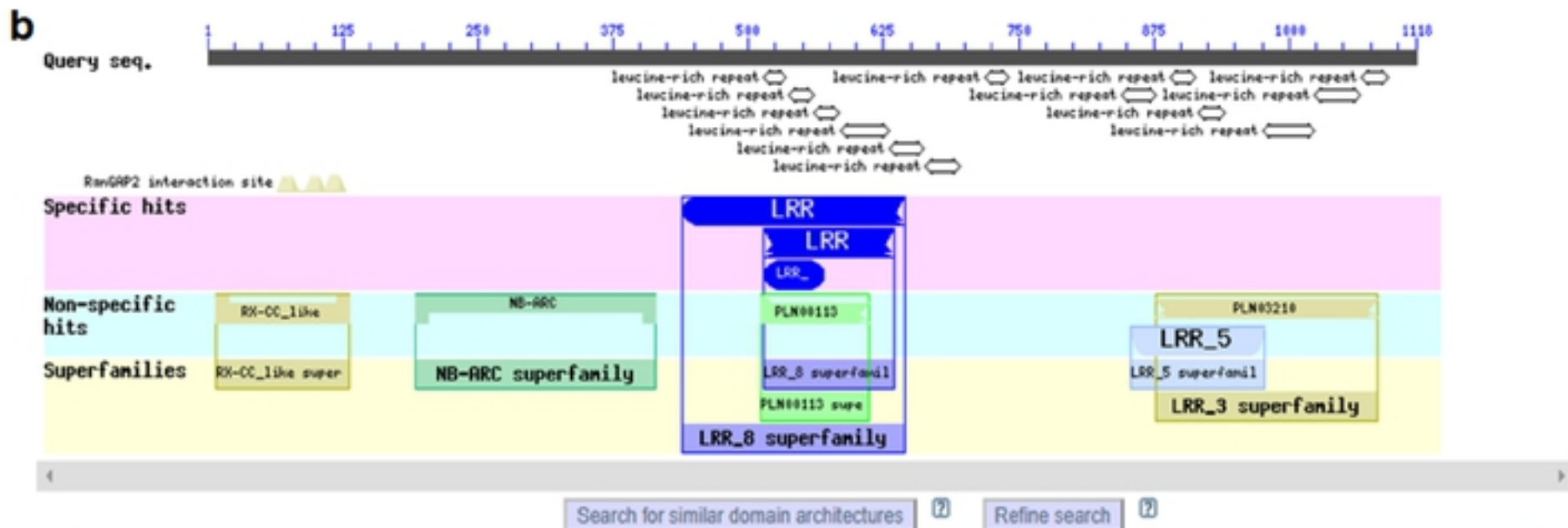714        doi:10.1073/pnas.1614862113

715   65.   Velasco R, Zharkikh A, Troggio M, Cartwright DA, Cestaro A, Pruss D, et al. A High Quality Draft

716        Consensus Sequence of the Genome of a Heterozygous Grapevine Variety. Dilkes B, editor. PLoS

717        One. 2007;2: e1326. doi:10.1371/journal.pone.0001326

718   66.   Sara M, Walter S, Paola C, Luigi M, Luigi F, Raffalella EM. Solanaceae Evolutionary Dynamics of

719        the &amp;lt;i&amp;gt;I&amp;lt;/i&amp;gt;2-NBS Domain. Am J Plant Sci. 2012;03: 283–294.

720        doi:10.4236/ajps.2012.32034

721   67.   Felsenstein J. Confidence-Limits on Phylogenies - an Approach Using the Bootstrap. Evolution (N

722        Y). 1985; doi:Doi 10.2307/2408678

723   68.   Zuckerkandl E, Pauling L. Evolutionary divergence and convergence in proteins. Evol genes proteins.

724        1965; 97–166. doi:10.1209/epl/i1998-00224-x

725   69.   Nei M, Kumar S. Molecular evolution and phylogenetics. 1st ed. Oxford: Oxford University Press;

726        2000.

727   70.   Saitou N, Nei M. The neighbor-joining method: a new method for reconstructing phylogenetic trees.

728        Mol Biol Evol. Oxford University Press; 1987;4: 406–25. Available:

729        http://www.ncbi.nlm.nih.gov/pubmed/3447015

a

| | | | |
|---|---|---|---|
| **Name** | **Accession** | **Description** | **Interval** | **E-value** |

| Name | Accession | Description | Interval | E-value |
|---|---|---|---|---|
| [+] NB-ARC | pfam00931 | NB-ARC domain; | 192-455 | 8.82e-46 |
| [+] PLN00113 | PLN00113 | leucine-rich repeat receptor-like protein kinase; Provisional | 507-655 | 1.03e-06 |
| [+] RX-CC_like | cd14798 | Coiled-coil domain of the potato virux X resistance protein and similar proteins; The potato ... | 2-132 | 3.45e-06 |
| [+] LRR | COG4886 | Leucine-rich repeat (LRR) protein [Transcription]; | 545-688 | 3.26e-04 |
| [+] CoaE | COG0237 | Dephospho-CoA kinase [Coenzyme transport and metabolism]; | 192-283 | 1.13e-03 |
| [+] PLN03210 | PLN03210 | Resistant to P. syringae 6; Provisional | 880-1117 | 2.98e-03 |

**b**



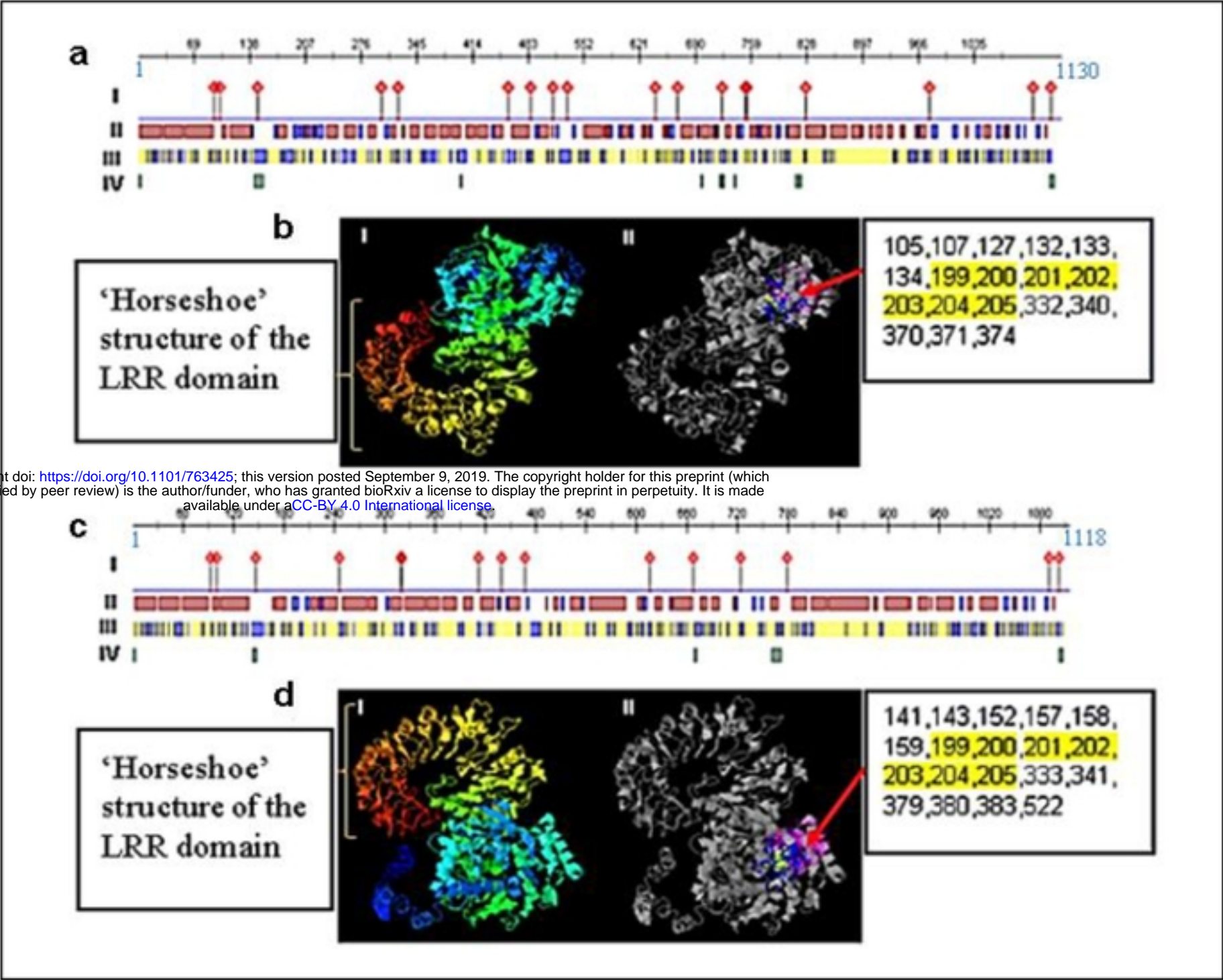| Name | Accession | Description | Interval | E-value |
|---|---|---|---|---|
| [+] NB-ARC | pfam00931 | NB-ARC domain; | 192-414 | 3.15e-37 |
| [+] PLN03210 | PLN03210 | Resistant to P. syringae 6; Provisional | 877-1082 | 1.04e-06 |
| [+] PLN00113 | PLN00113 | leucine-rich repeat receptor-like protein kinase; Provisional | 512-612 | 1.26e-05 |
| [+] LRR | COG4886 | Leucine-rich repeat (LRR) protein [Transcription]; | 439-645 | 1.37e-05 |
| [+] RX-CC_like | cd14798 | Coiled-coil domain of the potato virux X resistance protein and similar proteins; The potato ... | 8-132 | 2.00e-05 |
| [+] LRR | COG4886 | Leucine-rich repeat (LRR) protein [Transcription]; | 515-635 | 1.29e-04 |
| [+] LRR_8 | pfam13855 | Leucine rich repeat; | 515-571 | 2.14e-04 |
| [+] LRR_5 | pfam13306 | Leucine rich repeats (6 copies); This family includes a number of leucine rich repeats. This ... | 854-977 | 2.12e-03 |

```
gene5_aa    MADAAVSATVKAVLGTVISIAADRVGMVLGVKAELERLGKTTATIQGFLADADEKMHSQG    60
gene11_aa   MADTVISATVEVVLGTVISIAADRIGMARGVKAELERLSKTAAMMQGFLADCDEKMHTRG    60
            ***:..:****:. *********:.**.  *********.**:* :******.*****::*

gene5_aa    VRGWLKELEDEVFKADNVLDELHYHNLRQEVKYRNQPMKKKVCFFFSFFNAIGFSSSLAS   120
gene11_aa   VREWLKQLEDEVFKADNVLDELNYNNLRWDVKYRNQPMKKKVCFFFSFFSSIGFSSSLAS   120
            ** ***.:*****************.*:.***  .:***********************. .:*********

gene5_aa    KIRDINTNLERINQQANELGLVRKHQKEADAAGATASRQTDSIVVPNVVGRAVDESKIVE   180
gene11_aa   KIRDINTNLERINRQANELGLVRKHQKEANATGATTSRPTDSIVVPNVVGRAGDESKIVE   180
            *************:.********.********:.*.:***.**  ***************  *******

gene5_aa    MLLTPSERVVSVIPITGMGGLGKTTLAKSVYNNTKIVENFGIKSWVCVAREIKIVELFKL   240
gene11_aa   MLLTPSEKVVSVIPITGMGGLGKTTLAKSVYNNTKIDENFGIKSWVCVAREIKIVELFKL   240
            *******:*****************************  *********************

gene5_aa    ILESLPGTKVEVDGREAIVQEIRRKLGEKRFLLVLDDVWNRQWGLWNDFFTTLLGLSTTK   300
gene11_aa   ILESLTRTKVEVDGRDAIVQEIRGKLGEKRFLLVLDDVWNCEQEFWSDFFTTLLGLSTTK   300
            *****  ********:******* ****************  :  :*.************

gene5_aa    GSWCILTTRLEPVANAVPRHLQMND-PYFLGKLSDDACWSMLKEQVIAGEEVPQELEAIQ   359
gene11_aa   GSWCILTTRLQPVANAVPRHLQMNDGPYFLGKLSDDACWSILEKLVVAGEEVPNELEALK   360
            **********:***************  ************:*:: *:******:****::
```

```
gene5_aa    EQILRRCDGLPLAASLIGGLLLNNRKEKWHCIVQESLLNEDQGEIDQILKVSFDHLSPPS    419
gene11_aa   KQILKKCDGLPLAAKLIG------------------------------------------    378
            :.***::*******.***

gene5_aa    VKKCFAYCSIFPQDTKLGEDELIELWVAEGFVLPDRENTGMIEERGGEYLRILLQSSLLE    479
gene11_aa   VKKCFAYCSIFPQDTELGEDELIEHWVAEGFVLPDQKNTRMMEETGGEYLRILLQNSLLE    438
            **************** :******* ********** :** *:** **********.****

gene5_aa    KVADEGRTYYKMHDLVHDFAKSVLNPKSSSQDRYLALHSYEEMAENVRRNKAASIRSLFL    539
gene11_aa   KVQDKLRTYYKMHDLVHDFAKSILNPESSNQDRYLALNSSEGLVEKTTMTIPASIRTLFL    498
            ** *.* ****************:***.**.******.* * :. *:.  ****:***

gene5_aa    HSGGGISADMNMLSRFKHLHVLKLSGYDVVFLPSSIGKLLRLRLLDISSSGITSLPESLC    599
gene11_aa   HLEDGISAG--MLLRFKYLHVLRLSGNDVVFLPSSIGKLLHLRLLDISSSRIKSLPESLC    556
            *  .****.  ** ***:****:***.***********.**********  *.*******

gene5_aa    KLYNLQTLTIGGYALEGGFPKRMSDLISLRHLNYYHDDTEFKMLVQIGRLTCLQTLEFFN    659
gene11_aa   KLYNLQTLTIRNNALGEGFPKRMNDLISLRHLNYYHRAKFKMPMQMGQLTCLQTLKFFN    616
            **********  . **  ******.***************  :.***  .:.*:*.*********.***

gene5_aa    VSQEKGCGIEELGTLKYLKGSLEIRNLGLVKGKEAAKQAKLFEKPNLSRLVFKWESNL-S    718
gene11_aa   VSQEKGCGIEELGTLKYLRGSLEIRNLGLVEGKDAAKQAKLFEKPNLSRLRLDFRRKRGH    676
            ******************:***********:**:****************  :.:.*
```

```
gene5_aa    QKSDNRDEDVLEGLQPHPKLEKLKIGSFMGNKFPQWLINLPKLVVLRIEDCGRCSELPAL      778
gene11_aa   RKSDNCDEDVFEGLQPHPNLQKLEIRYFMGTKFSQWLINLPKLVELWIEDCKRCSELPSL      736
            *:**** ****:.*:****** *  *:: :::** ** ********** * **** **:: *

gene5_aa    GQLPSLKRLCLKRLENIRYVGDEFYGITTNE-----GSSRASGSSARRRKFFPALEKLKV      833
gene11_aa   GQLPSLKRLYLNKLENIRSIGDEFYGITTNEEGEEKGRSRASGSSTRRRKFFPALEELRV      796
            ********* *:.*****  :************      *  *******:************.*:*

gene5_aa    AFMENLAEWKDADQVRSTIGE--ADVFPMLRNFHIQSCPQLTALPCSCKILDVENCRNIT      891
gene11_aa   AYMKNLVEWKDADQVRSTIAEEAADVFPMLMDLSIQHCPQLTTLPCSCKILDVQYCRNLT      856
            *:*:**.*************.*    *******  :: ** *****:*********** : ***:*

gene5_aa    SIKTSYGTACVERLGIYSCDNLRELPVDVFGLSLQCLTISCCPRLISLGVNGKKCPLRC-      950
gene11_aa   SIKTGYGTASVEKLKIGCCNNLRELPEDVFGSSLQRLSIESCPRLISLGVNGKKCPLPCL      916
            ****.****.**:* *  .*:*****  ***  *** *** *:*.. **************** *

gene5_aa    ----------------------RSLRSVWVVSCPNLVSFSLNLQETPSLEEFVLDDCPKL      988
gene11_aa   ERLSIQYCYGLTTISDKMFESCQSLRSLSVECCPNLVSFSLNLQETPSLEDFALLNCPKL      976
                                  :****:  * .******************.*.* :****

gene5_aa    IPHNFKGFAFATSLRKLAIGPFSSDDSSIDDFDWSGLRSASTLRELYLQGLPRSKSLPHQ     1048
gene11_aa   IPHRFNGFAFATSLRNLWIGPFSSDDSSIDGFDWSGLRSASTLCKVHLEGLCHSDSLPHQ     1036
            ***.*:*********.* * ************* ************* :::*:** :*.*****

gene5_aa    LQYLATLTSLSLADFGGIEVLPDWIGNLVSLETLELSDCRKLQSLPSEAAMRRLTKLTHV     1108
gene11_aa   LQYLTTLTSLNLKNFGRIEVLPDWIGNLVSLETLQLSNCEKLRCLPSEAAMRRLTKLTSV     1096
            ****:*****.*  :.** ***************:**:*.**: .*************** *

gene5_aa    QVDGCPLLRQRYSPQRGIYLEE        1130
gene11_aa   EVRRCPLLRQRYTPQRGIYLEE        1118
            :* *******:.*********
```

a

‘Horseshoe’ structure of the LRR domain

105,107,127,132,133, 134,199,200,201,202, 203,204,205,332,340, 370,371,374

c

‘Horseshoe’ structure of the LRR domain

141,143,152,157,158, 159,199,200,201,202, 203,204,205,333,341, 379,380,383,522

Figure

Figure